

may also be applied to structures where semiconductor layer 106 is p-type and semiconductor strips 130 are n-type. During the epitaxy of p-type semiconductor strips 130, a p-type APT impurity (e.g., phosphorus or arsenic) may be in-situ doped with the proceeding or the epitaxy. Furthermore, p-type semiconductor strips 130 may be compressively strained as described above.

The epitaxy of p-type semiconductor strips 130 may overgrow trenches 128 and extend into opening 126. In FIG. 12, a high-selective CMP may be used to remove such overgrown portions and planarize top surfaces of p-type semiconductor strips 130. For example, the CMP process may comprise using a chemical slurry that selectively removes the material of p-type semiconductor strips 130 (e.g., SiGe or Ge) at a higher rate than barrier layer 122 (e.g., an oxide) and/or STI regions 120. In some embodiments, the slurry may etch p-type semiconductor strips 130 at about five times (or even greater) a rate than barrier layer 122 and/or STI regions 120. After CMP, top surfaces of p-type semiconductor strips 130 and STI regions 120 in opening 126 may be substantially level. Furthermore, the selective CMP process may further recess or even remove barrier layer 122. In embodiments where barrier layer 122 is removed during the selective CMP, pad layers 108 and/or 110 may protect n-type semiconductor strips 106 from damage.

Exposed portions of STI regions 120 may also be recessed (e.g., distance D1 may be increased). For example, in FIG. 11, a total dimension of distance D1 as a result of various etching, pre-cleaning, and/or CMP processes may be about 20 nm or even more. As a further result of the various etching, pre-cleaning, and/or CMP processes used to form p-type semiconductor strips 130, top surfaces of p-type semiconductor strips 130 and n-type semiconductor strips 106 may not be substantially level. For example, a top surface of p-type semiconductor strips 130 may be recessed from a top surface of n-type semiconductor strips 106 by a distance D3. In some embodiments, D3 may be about 12 nm to about 22 nm, for example. Such a height variation between n-type semiconductor strips 106 and p-type semiconductor strips 130 may result in various manufacturing and/or device defects. Thus, additional processing may be used to further level top surfaces of n-type semiconductor strips 106 and p-type semiconductor strips 130.

FIGS. 13 through 15 illustrate further process steps to level top surfaces of n-type semiconductor strips 106 and p-type semiconductor strips 130. Referring to FIG. 13, a sacrificial layer 132 is blanket deposited over a top surface of wafer 100. In some embodiments, sacrificial layer 132 comprises a plasma-enhanced oxide (PEOX) deposited using any suitable method, such as CVD, PVD, ALD, and the like. Sacrificial layer 132 may be used to planarize a top surface of wafer 100 and to mitigate the height variation between n-type semiconductor strips 106 and p-type semiconductor strips 130. In some embodiments, sacrificial layer 132 is sufficiently thick to fill (and even overflow) opening 126. For example, sacrificial layer 132 may have a thickness D4 of about 200 nm to about 300 nm. Other dimensions may also be used depending on wafer configuration.

In FIG. 14, a planarization may be performed to remove barrier layer 122 and portions of sacrificial layer 132 extending over pad layer 110. For example, an end-point CMP process may be used to planarize a top surface of wafer 100. During the end-point CMP process, pad layer 110 (e.g., a nitride) may be used as an end point detection layer. During

the CMP, sacrificial layer 132 may be used to resist CMP and prevent damage to underlying p-type semiconductor strips 130.

Next, in FIG. 15, a further planarization may be performed to remove pad layer 110, pad layer 108, and sacrificial layer 132. For example, a non-selective CMP process may be used to further planarize wafer 100. In some embodiments, the non-selective CMP process may comprise using a chemical slurry which etches pad layer 110 (e.g., a nitride), pad layer 108/sacrificial layer 132 (e.g., oxides), and n-type semiconductor strips 106/p-type semiconductor strips 130 at a rate of about 1.1 to about 1 to about 0.5, for example.

After planarization, top surfaces of STI regions 120, n-type semiconductor strips 106, and p-type semiconductor strips 130 may be substantially level. It has been observed that by using the various processing steps described above with respect to FIGS. 1 through 15, improved fin height uniformity may be achieved. For example, in fins formed using the above described process steps, a height difference between top surfaces of n-type semiconductor strips 106 and p-type semiconductor strips 130 may be about 3 nm to about 5 nm, while fins formed using other processing methods may exhibit a height difference be about 12 nm to about 22 nm or even more. Furthermore, the use of various barrier layers/sacrificial layers allows for a relatively small amount of non-selective planarization techniques (e.g., non-selective CMP) while still achieving a suitably planar top surface. Thus, the total height of n-type semiconductor strips 106 and p-type semiconductor strips 130 may be maintained at a desired level (e.g., at a desired channel height), which allows for improved device performance and/or reliability in resulting finFETs. For example, in some embodiments, a height of n-type semiconductor strips 106 and p-type semiconductor strips 130 (e.g., distance D5) may be about 40 nm. Other dimensions may also be used depending on wafer configuration.

After the formation of fins 116 and 116' (e.g., comprising semiconductor strips of different types), additional processing may be performed to create finFETs in wafer 100. For example, in FIG. 16, STI regions 120 are recessed, so that top portions of semiconductor strips 106 and 130 are higher than the top surfaces of STI regions 120. The recessing of STI regions 120 may include a chemical etch process, for example, using ammonia (NH<sub>3</sub>) in combination with hydrofluoric acid (HF) or nitrogen trifluoride (NF<sub>3</sub>) as reaction solutions either with or without plasma. When HF is used as the reaction solution, a dilution ratio of HF may be between about 1:50 to about 1:100.

Channel regions of two different types (e.g., corresponding to n-type semiconductor strips 106 and p-type semiconductor strips 130) are thus formed in fins 116/116'. In the completed finFET structures, one or more gate stacks wrap around and covers sidewalls of such channel regions (see FIG. 17). The resulting structure is shown in FIG. 17. For example, referring to FIG. 17, gate stacks are formed on the top surface and the sidewalls of semiconductor strips 106 and 130. Such gate stacks may include a gate dielectric 152 and a gate electrode 154.

In accordance with some embodiments, gate dielectric 152 includes silicon oxide, silicon nitride, or multilayers thereof. In alternative embodiments, gate dielectric 152 includes a high-k dielectric material. In such embodiments, gate dielectric 50 may have a k value greater than about 7.0, and may include a metal oxide or a silicate of hafnium (Hf), aluminum (Al), zirconium (Zr), lanthanum (La), magnesium (Mg), barium (Ba), titanium (Ti), lead (Pb), and combina-